

ST243 Homework 3

Sicheng Zhu: 915640461

Ruoyan Yin: 916666619

Question 2

Part a

$$p(Z_i = j) = \pi_j$$

Part b

$$\begin{aligned} p(Z_i = j|x_i) &= \frac{f(x_i|Z_i = j)p(Z_i = j)}{\sum_{l=1}^k f(x_i|Z_i = l)p(Z_i = l)} \\ &= \frac{\frac{\pi_j}{(2\pi)^{d/2}|\Sigma_j|^{1/2}} \exp\{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1}(x_i - \mu_j)\}}{\sum_{l=1}^k \frac{\pi_l}{(2\pi)^{d/2}|\Sigma_l|^{1/2}} \exp\{-\frac{1}{2}(x_i - \mu_l)^T \Sigma_l^{-1}(x_i - \mu_l)\}} \\ &= \frac{\frac{\pi_j}{|\Sigma_j|^{1/2}} \exp\{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1}(x_i - \mu_j)\}}{\sum_{l=1}^k \frac{\pi_l}{|\Sigma_l|^{1/2}} \exp\{-\frac{1}{2}(x_i - \mu_l)^T \Sigma_l^{-1}(x_i - \mu_l)\}} \end{aligned}$$

Part c

$$F_{ij} = p_\theta(Z_i = j|x_j)$$

By Jensen's Inequality,

$$\begin{aligned} l(\theta) &= \sum_{i=1}^n \log \left[\sum_{j=1}^k F_{ij} \frac{p_\theta(x_i, Z_i = j)}{F_{ij}} \right] = \sum_{i=1}^n \log E \left[\frac{p_\theta(x_i, Z_i = j)}{q_i(Z_i)} \right] \\ &\geq \sum_{i=1}^n E \left[\log \frac{p_\theta(x_i, Z_i = j)}{q_i(Z_i)} \right] = \sum_{i=1}^n \sum_{j=1}^k F_{ij} \log \frac{p_\theta(x_i, Z_i = j)}{F_{ij}} \end{aligned}$$

Part d

In Jensen's Inequality, $E_\theta(\log X) = \log(E_\theta X) \Leftrightarrow X$ is constant to θ .

$$\begin{aligned} l(\theta') &= \sum_{i=1}^n \log \left[\sum_{j=1}^k F_{ij} \frac{p_{\theta'}(x_i, Z_i = j)}{F_{ij}} \right] \\ &= \sum_{i=1}^n \log \left[\sum_{j=1}^k p_{\theta'}(Z_i = j|x_i) \frac{p_{\theta'}(x_i, Z_i = j)}{p_{\theta'}(Z_i = j|x_i)} \right] \\ &= \sum_{i=1}^n \log \left[\sum_{j=1}^k p_{\theta'}(Z_i = j|x_i) f_{\theta'}(x_i) \right] = \sum_{i=1}^n \log E_{Z_i} [f_{\theta'}(x_i)] \\ &= \sum_{i=1}^n E_{Z_i} [\log f_{\theta'}(x_i)] \sum_{j=1}^k p_{\theta'}(Z_i = j|x_i) \log \frac{p_{\theta'}(x_i, Z_i = j)}{p_{\theta'}(Z_i = j|x_i)} \end{aligned}$$

Part e

$$\begin{aligned}
 Q(\theta, \theta^{(t)}) &= \sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \log \pi_j N(x_i, \mu_j, \Sigma_j) \\
 &= \sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \log \pi_j + \sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \log N(x_i, \mu_j, \Sigma_j) \\
 \pi^{(t+1)} &= \operatorname{argmax}_{\pi: \sum \pi_j = 1} \left(\sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \log \pi_j \right) = \operatorname{argmax}_{\pi: \sum \pi_j = 1} \sum_{i=1}^n \log \pi_j \left(\sum_{j=1}^k F_{ij}^{(t)} \right)
 \end{aligned}$$

$$\text{where } F_{ij}^{(t)} = p_{\theta^{(t)}}(Z_i = j | x_i) = \frac{\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{j'=1}^k \pi_{j'}^{(t)} N(x_i, \mu_{j'}^{(t)}, \Sigma_{j'}^{(t)})}$$

Calculating $\pi_j^{(t+1)}$:

$$\begin{aligned}
 \phi(\pi, \lambda) &= \sum_{j=1}^k \log \pi_j \left(\sum_{i=1}^n F_{ij}^{(t)} \right) + \left(\sum_{j=1}^k \pi_j - 1 \right) \\
 \begin{cases} \frac{\partial \phi}{\partial \pi_j} = \frac{1}{\pi_j} \sum_{i=1}^n F_{ij}^{(t)} + \lambda = 0, \forall j = 1, \dots, k \\ \frac{\partial \phi}{\partial \lambda} = \sum_{j=1}^k \pi_j - 1 = 0 \end{cases} &\Rightarrow \begin{cases} \lambda = - \sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \\ \pi_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)}}{\sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)}} = \frac{\sum_{i=1}^n F_{ij}^{(t)}}{n} \end{cases}
 \end{aligned}$$

Calculating $\mu_j^{(t+1)}$ and $\Sigma_j^{(t+1)}$:

$$\begin{aligned}
 (\mu^{(t+1)}, \Sigma^{(t+1)}) &= \sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \log N(x_i, \mu_j, \Sigma_j) \\
 &= \operatorname{argmax}_{\mu, \Sigma} \sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \log \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} e^{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)} \\
 &= \operatorname{argmax}_{\mu, \Sigma} \sum_{i=1}^n \sum_{j=1}^k F_{ij}^{(t)} \left(-\frac{1}{2} (x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j) - \frac{1}{2} \log |\Sigma_j| \right) \\
 &= \operatorname{argmax}_{\mu, \Sigma} \sum_{i=1}^n \sum_{j=1}^k -\frac{1}{2} F_{ij}^{(t)} (\log |\Sigma_j| + x_i^T \Sigma_j^{-1} x_i - \mu_j^T \Sigma_j^{-1} x_i - x_i^T \Sigma_j^{-1} \mu_j \\
 &\quad + \mu_j^T \Sigma_j^{-1} \mu_j)
 \end{aligned}$$

$$\frac{\partial Q}{\partial \mu_j} = \sum_{i=1}^n -\frac{1}{2} F_{ij}^{(t)} \left(-(\Sigma_j^{-1} x_i)^T - x_i^T \Sigma_j^{-1} + 2\mu_j^T \Sigma_j^{-1} \right) = \sum_{i=1}^n F_{ij}^{(t)} (x_i^T \Sigma_j^{-1} - \mu_j^T \Sigma_j^{-1}) = 0$$

$$\sum_{i=1}^n F_{ij}^{(t)} x_i^T \Sigma_j^{-1} = \sum_{i=1}^n F_{ij}^{(t)} \mu_j^T \Sigma_j^{-1}$$

$$\Sigma_j^{-1} \sum_{i=1}^n F_{ij}^{(t)} x_i^T = \Sigma_j^{-1} \sum_{i=1}^n F_{ij}^{(t)} \mu_j^T$$

$$\Rightarrow \mu_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)} x_i}{\sum_{i=1}^n F_{ij}^{(t)}}$$

$$\begin{aligned} \Sigma_j^{(t+1)} &= \operatorname{argmax}_{\Sigma} \sum_{i=1}^n -\frac{1}{2} F_{ij}^{(t)} \left((x_i - \mu_j^{(t+1)})^T \Sigma_j^{-1} (x_i - \mu_j^{(t+1)}) + \log |\Sigma_j| \right) \\ &= \operatorname{argmax}_{\Sigma} \left[-\frac{1}{2} \sum_{i=1}^n \operatorname{tr} \left(F_{ij}^{(t)} \Sigma_j^{-1} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T \right) - \frac{1}{2} \log |\Sigma_j| \sum_{i=1}^n F_{ij}^{(t)} \right] \\ &= \operatorname{argmax}_{\Sigma} \left[-\frac{1}{2} \operatorname{tr} \left(\Sigma_j^{-1} \sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T \right) - \frac{1}{2} \log |\Sigma_j| \sum_{i=1}^n F_{ij}^{(t)} \right] \end{aligned}$$

$$\frac{\partial Q}{\partial \Sigma_j} = -\frac{1}{2} \left[-(\Sigma_j^{-1})^T \sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T (\Sigma_j^{-1})^T + \Sigma_j^{-1} \sum_{i=1}^n F_{ij}^{(t)} \right] = 0$$

As Σ_j is diagonal matrix, Σ_j^{-1} is also diagonal matrix, $(\Sigma_j^{-1})^T = \Sigma_j^{-1}$

$$\Sigma_j^{-1} \sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T \Sigma_j^{-1} = \sum_{i=1}^n F_{ij}^{(t)} \Sigma_j^{-1}$$

$$\Sigma_j \Sigma_j^{-1} \sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T \Sigma_j^{-1} \Sigma_j = \Sigma_j \sum_{i=1}^n F_{ij}^{(t)} \Sigma_j^{-1} \Sigma_j$$

$$\sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T = \Sigma_j \sum_{i=1}^n F_{ij}^{(t)}$$

$$\Rightarrow \Sigma_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T}{\sum_{i=1}^n F_{ij}^{(t)}}$$

$$\begin{cases} \pi_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)}}{n} \\ \mu_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)} x_i}{\sum_{i=1}^n F_{ij}^{(t)}} \\ \Sigma_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T}{\sum_{i=1}^n F_{ij}^{(t)}} \end{cases}$$

In a mixture of spherical Gaussians model, $\Sigma_j = \sigma_j^2 I_d$, $|\Sigma_j| = \sigma_j^{2d}$, $\Sigma_j^{-1} = \frac{1}{\sigma_j^2} I_d$

$$\begin{aligned} N(x_i, \mu_j, \Sigma_j) &= \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} e^{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)} \\ &= \frac{1}{(2\pi)^{d/2} \sigma_j^d} e^{-\frac{1}{2}(x_i - \mu_j)^T \frac{1}{\sigma_j^2} I_d (x_i - \mu_j)} = \frac{1}{(2\pi)^{d/2} \sigma_j^d} e^{-\frac{1}{2\sigma_j^2} (x_i - \mu_j)^T (x_i - \mu_j)} \end{aligned}$$

So,

$$\begin{aligned} F_{ij}^{(t)} &= \frac{\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{j'=1}^k \pi_{j'}^{(t)} N(x_i, \mu_{j'}^{(t)}, \Sigma_{j'}^{(t)})} \\ &= \frac{\pi_j^{(t)} \frac{1}{(2\pi)^{d/2} \sigma_j^{(t)d}} \exp\{-\frac{1}{2\sigma_j^{(t)2}} (x_i - \mu_j^{(t)})^T (x_i - \mu_j^{(t)})\}}{\sum_{j'=1}^k \pi_{j'}^{(t)} \frac{1}{(2\pi)^{d/2} \sigma_{j'}^{(t)d}} \exp\{-\frac{1}{2\sigma_{j'}^{(t)2}} (x_i - \mu_{j'}^{(t)})^T (x_i - \mu_{j'}^{(t)})\}} \\ &= \frac{\pi_j^{(t)} \frac{1}{\sigma_j^{(t)d}} \exp\{-\frac{1}{2\sigma_j^{(t)2}} (x_i - \mu_j^{(t)})^T (x_i - \mu_j^{(t)})\}}{\sum_{j'=1}^k \pi_{j'}^{(t)} \frac{1}{\sigma_{j'}^{(t)d}} \exp\{-\frac{1}{2\sigma_{j'}^{(t)2}} (x_i - \mu_{j'}^{(t)})^T (x_i - \mu_{j'}^{(t)})\}} \end{aligned}$$

$$\begin{aligned} \sigma_j^{(t+1)} &= \underset{\sigma}{\operatorname{argmax}} \sum_{i=1}^n -\frac{1}{2} F_{ij}^{(t)} \left(\frac{1}{\sigma_j^2} (x_i - \mu_j^{(t+1)})^T (x_i - \mu_j^{(t+1)}) + \log \sigma_j^{2d} \right) \\ &= \underset{\sigma}{\operatorname{argmax}} \sum_{i=1}^n -\frac{1}{2} F_{ij}^{(t)} \frac{1}{\sigma_j^2} (x_i - \mu_j^{(t+1)})^T (x_i - \mu_j^{(t+1)}) - d F_{ij}^{(t)} \log \sigma_j \end{aligned}$$

$$\frac{\partial Q}{\partial \sigma_j} = \frac{1}{\sigma_j^3} \sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)})^T (x_i - \mu_j^{(t+1)}) - \frac{d}{\sigma_j} \sum_{i=1}^n F_{ij}^{(t)} = 0$$

$$\begin{cases} \pi_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)}}{n} \\ \mu_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)} x_i}{\sum_{i=1}^n F_{ij}^{(t)}} \\ \sigma_j^{2(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)} (x_i - \mu_j^{(t+1)})^T (x_i - \mu_j^{(t+1)})}{d \sum_{i=1}^n F_{ij}^{(t)}} \end{cases},$$

$$\text{where } F_{ij}^{(t)} = \frac{\pi_j^{(t)} \frac{1}{\sigma_j^{(t)d}} \exp\{-\frac{1}{2\sigma_j^{(t)2}}(x_i - \mu_j^{(t)})^T (x_i - \mu_j^{(t)})\}}{\sum_{j'=1}^k \pi_{j'}^{(t)} \frac{1}{\sigma_{j'}^{(t)d}} \exp\{-\frac{1}{2\sigma_{j'}^{(t)2}}(x_i - \mu_{j'}^{(t)})^T (x_i - \mu_{j'}^{(t)})\}}$$

Part f

In a mixture of diagonal Gaussians model, $\Sigma_j = \text{diag}(\sigma_{j1}^2, \sigma_{j2}^2, \dots, \sigma_{jd}^2)$,

$$|\Sigma_j| = \prod_{l=1}^d \sigma_{jl}^2, \Sigma_j^{-1} = \text{diag}\left(\frac{1}{\sigma_{j1}^2}, \frac{1}{\sigma_{j2}^2}, \dots, \frac{1}{\sigma_{jd}^2}\right)$$

Let $\sigma_j^2 = [\sigma_{j1}^2, \sigma_{j2}^2, \dots, \sigma_{jd}^2]$, and $\frac{1}{\sigma_j^2} = [\frac{1}{\sigma_{j1}^2}, \frac{1}{\sigma_{j2}^2}, \dots, \frac{1}{\sigma_{jd}^2}]$

$$\begin{aligned} N(x_i, \mu_j, \Sigma_j) &= \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} e^{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)} \\ &= \frac{1}{(2\pi)^{d/2} \prod_{l=1}^d \sigma_{jl}} e^{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)} \\ &= \frac{1}{(2\pi)^{d/2} \prod_{l=1}^d \sigma_{jl}} e^{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)} \\ &= \frac{1}{(2\pi)^{d/2} \prod_{l=1}^d \sigma_{jl}} \exp\left\{\sum_{l=1}^d -\frac{1}{2\sigma_{jl}^2} (x_{il} - \mu_{jl})^T (x_{il} - \mu_{jl})\right\} \end{aligned}$$

So,

$$\begin{aligned} F_{ij}^{(t)} &= \frac{\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{j'=1}^k \pi_{j'}^{(t)} N(x_i, \mu_{j'}^{(t)}, \Sigma_{j'}^{(t)})} \\ &= \frac{\pi_j^{(t)} \frac{1}{(2\pi)^{d/2} \prod_{l=1}^d \sigma_{jl}^{(t)}} \exp\{\sum_{l=1}^d -\frac{1}{2\sigma_{jl}^{(t)2}} (x_{il} - \mu_{jl})^2\}}{\sum_{j'=1}^k \pi_{j'}^{(t)} \frac{1}{(2\pi)^{d/2} \prod_{l=1}^d \sigma_{j'l}^{(t)}} \exp\{\sum_{l=1}^d -\frac{1}{2\sigma_{j'l}^{(t)2}} (x_{il} - \mu_{j'l})^2\}} \end{aligned}$$

$$= \frac{\pi_j^{(t)} \frac{1}{\prod_{l=1}^k \sigma_{jl}^{(t)}} \exp\{\sum_{l=1}^d -\frac{1}{2\sigma_{jl}^{(t)2}} (x_{il} - \mu_{jl})^2\}}{\sum_{j'=1}^k \pi_{j'}^{(t)} \frac{1}{\prod_{l=1}^k \sigma_{j'l}^{(t)}} \exp\{\sum_{l=1}^d -\frac{1}{2\sigma_{j'l}^{(t)2}} (x_{il} - \mu_{j'l})^2\}}$$

$$\begin{aligned} \sigma_{jl}^{(t+1)} &= \operatorname{argmax}_{\sigma_{jl}} \sum_{i=1}^n -\frac{1}{2\sigma_{jl}^2} F_{ij}^{(t)} \left((x_{il} - \mu_{jl}^{(t+1)})^2 + 2 \log \sigma_{jl} \right) \\ &= \operatorname{argmax}_{\sigma_{jl}} \left[-\frac{1}{2} \sum_{i=1}^n \frac{1}{\sigma_{jl}^2} F_{ij}^{(t)} (x_{il} - \mu_{jl}^{(t+1)})^2 - \frac{1}{2} \log \sigma_{jl} \sum_{i=1}^n F_{ij}^{(t)} \right] \end{aligned}$$

$$\frac{\partial Q}{\partial \sigma_{jl}} = \frac{1}{\sigma_{jl}^3} \sum_{i=1}^n F_{ij}^{(t)} (x_{il} - \mu_{jl}^{(t+1)})^2 - \frac{1}{2\sigma_{jl}} \sum_{i=1}^n F_{ij}^{(t)} = 0$$

$$\begin{cases} \pi_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)}}{n} \\ \mu_j^{(t+1)} = \frac{\sum_{i=1}^n F_{ij}^{(t)} x_i}{\sum_{i=1}^n F_{ij}^{(t)}} \\ \sigma_{jl}^{2(t+1)} = \frac{2 \sum_{i=1}^n F_{ij}^{(t)} (x_{il} - \mu_{jl}^{(t+1)})^2}{\sum_{i=1}^n F_{ij}^{(t)}}, l = 1, \dots, d \end{cases},$$

$$\text{where } F_{ij}^{(t)} = \frac{\pi_j^{(t)} \frac{1}{\prod_{l=1}^k \sigma_{jl}^{(t)}} \exp\{\sum_{l=1}^d -\frac{1}{2\sigma_{jl}^{(t)2}} (x_{il} - \mu_{jl})^2\}}{\sum_{j'=1}^k \pi_{j'}^{(t)} \frac{1}{\prod_{l=1}^k \sigma_{j'l}^{(t)}} \exp\{\sum_{l=1}^d -\frac{1}{2\sigma_{j'l}^{(t)2}} (x_{il} - \mu_{j'l})^2\}}$$

Question 3

(i)

The EM algorithm for mixture of spherical Gaussians is implemented. And because each x_i is an 1x196 row vector, the value of $f(x_i, \mu_j, \Sigma_j)$ could be extremely small. In order to get the precise value of $F_{ij}^{(t)}$, we apply the log-sum-exp trick as below:

$$\begin{aligned}\log F_{ij}^{(t)} &= \log \left(\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)}) \right) - \log \left(\sum_{j'=1}^k \pi_{j'}^{(t)} N(x_i, \mu_{j'}^{(t)}, \Sigma_{j'}^{(t)}) \right) \\ &= \log \left(\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)}) \right) - \log \left(\sum_{j'=1}^k \exp\{\log(\pi_{j'}^{(t)} N(x_i, \mu_{j'}^{(t)}, \Sigma_{j'}^{(t)}))\} \right) \\ &= \log \left(\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)}) \right) - [a + \log \sum_{j'=1}^k \exp\{\log(\pi_{j'}^{(t)} N(x_i, \mu_{j'}^{(t)}, \Sigma_{j'}^{(t)})) - a\}] \end{aligned}$$

where $a = \max_j \log \left(\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)}) \right)$,

$$\log \left(\pi_j^{(t)} N(x_i, \mu_j^{(t)}, \Sigma_j^{(t)}) \right) = \log \pi_j^{(t)} - d \log \sigma_j^{(t)} - \frac{1}{2\sigma_j^{(t)2}} (x_i - \mu_j^{(t)})^T (x_i - \mu_j^{(t)}).$$

Then, within each iteration, we calculate the matrix F, and compute

$\pi_j^{(t+1)}, \mu_j^{(t+1)}, \sigma_j^{(t+1)}$ using the result of **part e** of Question 2.

As we need the fractional change of the log-likelihood between each iteration to decide terminating or not, for the first iteration, in order to let

$\frac{\log\text{-likelihood}_{(t)} - \log\text{-likelihood}_{(t-1)}}{\log\text{-likelihood}_{(t-1)}}$ work well for $t = 1$, we set $\log\text{-likelihood}_{(0)} =$

1, to make sure the denominator is not 0.

We tried 3 initializations:

	I	II	III
$\pi^{(0)}$	[0.2 0.2 0.2 0.2 0.2]	[0.1 0.4 0.4 0.02 0.08]	[0.8 0.05 0.05 0.05 0.05]
$\mu_{jl}^{(0)}$	$\sim N(100, 10)$	$\sim N(50, 5)$	$\sim N(0, 1)$
$\sigma_j^{(0)}$	100	$\sim N(50, 10)$	$\sim N(200, 10)$

The iteration times and log-likelihoods results are as below:

	I	II	III
<i>Iteration times</i>	12	19	11
<i>log-likelihood</i>	-31827008	-31795742	-31827819

The second iteration got the best result in terms of maximum log-likelihood. The third iteration got the best result in terms of speed of convergence.

We identify the digits of clusters by the following method:

- (a) Assume the digit of the cluster to be i ;
- (b) Count the number of correctly classified observations t ;
- (c) Compare t of different i 's, and assign the cluster to be digit i with largest t .

We plot each μ_j and the digits showing in the below images are exactly the same as what we identify from the above methods :

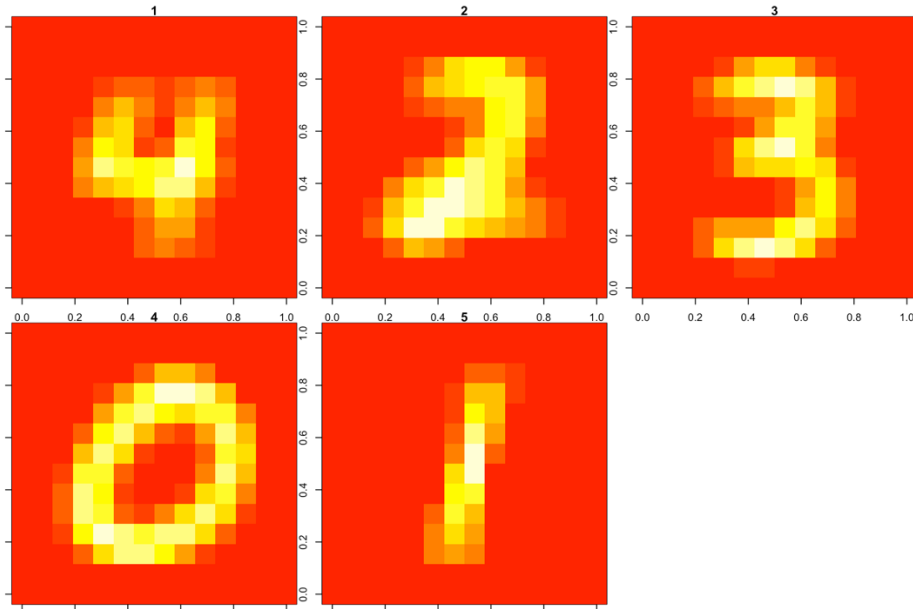


Figure 1 Initialization I

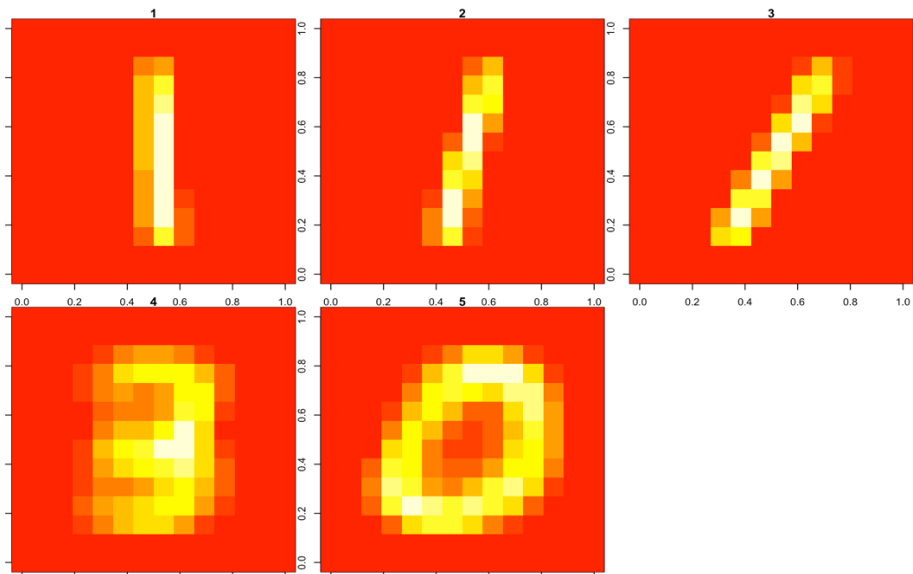


Figure 2 Initialization II

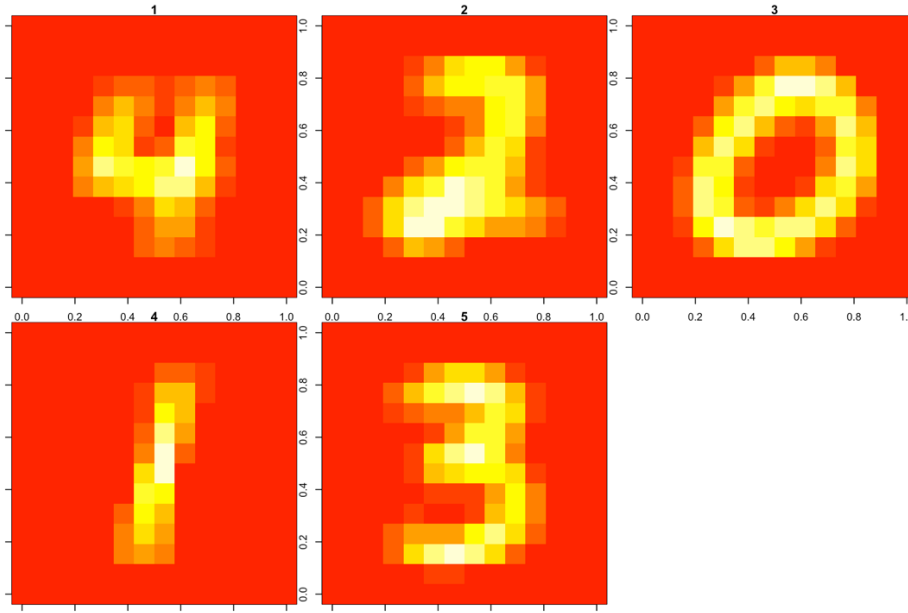


Figure 3 Initialization III

For initialization I and III, the clusters are successful, as five clusters refer to five digits. But for initialization II, 3 clusters are digit 1, and it misses digits 2 and 4. Although the log-likelihood from initialization II is minimal among all the three, the result of it is bad, and the errors should be very large such that we do not test over this model.

Using EM algorithms from the three initializations to fit the testing data, the errors are as below:

	I	III
<i>Errors</i>	0.1152	0.1136

The clusters coming from initialization III has the best result in terms of prediction error. Still, the error is significantly large, so a mixture models of Gaussian distributions may not be suitable for modeling the MNIST data.

(ii)

Then we applied EM algorithm for diagonal Gaussians. Still, we used log-exp trick to avoid the problem that the conditional probability and the log-likelihood can be underflow. Just like what we did in part (i), we still set the proportional change in log-likelihood as criterion to make the decision if the classification is good enough.

Since the outcome can be highly dependent on the initialized parameters, we tried three different combinations.

	I	II	III
π	(0.2 0.25 0.25 0.15 0.15)	(0.15 0.20 0.25 0.20 0.20)	(0.3 0.2 0.1 0.20 0.20)
μ	$\sim N(80, 2)$	$\sim N(100, 10)$	$\sim N(60, 5)$
σ	uniform(0.9, 1.1)	uniform(0.5, 1.3)	uniform(0.8, 1.5)

Chart below shows the iteration numbers taken by different combinations and the log-likelihood based on those initializations.

	I	II	III
<i>Iteration times</i>	7	11	11
log-likelihood	-21485205	-19067068	--19658301

Even the second combination get the best result with respect to the log-likelihood, it didn't successfully depart the data into 5 groups. And the third one got a better likelihood even it takes a few steps than the first.

Using EM algorithms from the three initializations to fit the testing data, the errors are as below:

	I	III
<i>Errors</i>	0.5606149	0.5197509

The clusters coming from initialization III has a relative better result in terms of prediction error compared to I. However, the error is significantly large, so a mixture models of Gaussian distributions may not be suitable for modeling the MNIST data.